

ADAPTED MULTIMODAL BERT WITH LAYER-WISE FUSION FOR SENTIMENT ANALYSIS

Odysseas S. Chlapanis¹ Georgios Paraskevopoulos^{1,2} Alexandros Potamianos¹

¹ National Technical University of Athens, Athens, Greece

² Institute for Language and Speech Processing, Athena Research Center, Athens, Greece

ABSTRACT

Multimodal learning pipelines have benefited from the success of pretrained language models. However, this comes at the cost of increased model parameters. In this work, we propose Adapted Multimodal BERT (AMB), a BERT-based architecture for multimodal tasks that uses a combination of adapter modules and intermediate fusion layers. The adapter adjusts the pretrained language model for the task at hand, while the fusion layers perform task-specific, layer-wise fusion of audio-visual information with textual BERT representations. During the adaptation process the pre-trained language model parameters remain frozen, allowing for fast, parameter-efficient training. In our ablations we see that this approach leads to efficient models, that can outperform their fine-tuned counterparts and are robust to input noise. Our experiments on sentiment analysis with CMU-MOSEI show that AMB outperforms the current state-of-the-art across metrics, with 3.4% relative reduction in the resulting error and 2.1% relative improvement in 7-class classification accuracy.

Index Terms— adapters, BERT, multimodal, fusion

1. INTRODUCTION

Over the past few years, we have witnessed impressive breakthroughs in the field of multimodal applications, due to the abundance of multimedia data and progress in core machine learning algorithms. This has set the scene for multimodal machine learning as one of the frontiers of applied AI research. For wide-spread adoption in the real-world, models that strike the correct balance between performance and parameter efficiency should be developed.

GPT [1] and BERT [2] were the first to establish the effectiveness of pre-training large scale language models on general tasks and then refining them for a specific task. Inspired by this approach, ViBERT [3] leveraged parallel multimodal data for pre-training a visual-language model. Other researchers [4, 5, 6, 7] have adopted a more flexible method: adapting a model pre-trained only on language for multimodal tasks.

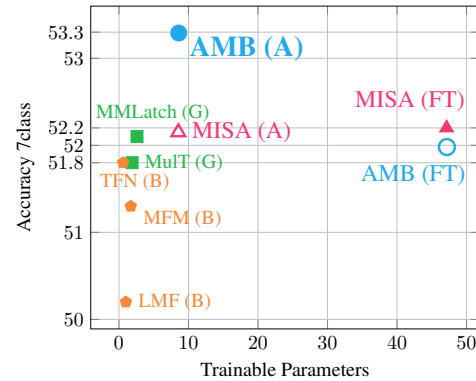


Fig. 1: 7-class accuracy with respect to number of trainable parameters for the best performing models in the literature. G stands for GloVe embeddings, A for adapters, B for frozen and FT for fine-tuned BERT embeddings. The proposed AMB with adapters achieves a good balance between trainable parameters and performance.

The standard method of transferring a pre-trained model to a downstream task is called fine-tuning, which involves updating the pre-trained weights with backpropagation. However this method incurs intensive data and computational costs, while some information is lost due to using only task-specific data for updating model parameters. This phenomenon is known as catastrophic forgetting [8]. To solve these issues, GPT-3 [9] proposed “prompt tuning”, an intuitive method to transfer a powerful pretrained model only with text interactions, called “prompts”, without any gradient updates. This idea was later extended, with many variations [10, 11], to make these prompts trainable, now called “soft prompts”. Houlsby et al. [12] proposed adapters, a down-projected feedforward network that updates the representations of each BERT layer. Frozen [4] applied these ideas in multimodal learning, by translating an image to a visual soft prompt that is prepended to the input of a standard language model, which keeps its original pre-trained weights unchanged (frozen). MAGMA [5] extended this by showing that the addition of adapter layers [12] in between the frozen language layers outperforms Frozen. Flamingo [13] scaled up and optimised this concept by introducing a flexible visual

encoder which can turn arbitrary sequences of images or even video frames to a fixed number of visual tokens.

Early applications of deep learning for multimodal sentiment analysis focused on the use of Recurrent Neural Networks (RNNs) [14, 15, 16] and Convolutional Neural Networks (CNNs) [17] aiming to model contextual information. The next innovation was the introduction of the attention model to create sophisticated fusion approaches [18, 19]. This naturally led to the incorporation of the transformer [20] as the central model for this task [21, 22]. Lately, large-scale pretrained language transformers, such as BERT [2], have become the norm because of consistent performance gains. ICCN [23] introduced Deep Canonical Correlation Analysis for jointly learning representations. Wang et al. [19] and later MAG-BERT [7] proposed shifting methods. MISA [6] produced modality invariant and modality specific representations in an effort to disentangle data relationships. More recently, many researchers turned their efforts towards intricate multimodal pre-training strategies, such as [24, 25]. Such methods are model-agnostic and should be studied separately for a fair comparison.

We present a simple neural architecture that adapts BERT representations for multimodal fusion which we call Adapted Multimodal BERT (AMB). Our approach extends concepts introduced by visual-language models [4, 5, 13] to include audio. The contributions of our work:

- AMB is evaluated on multimodal sentiment analysis with the CMU-MOSEI database to achieve new state-of-the-art results, regardless of being lightweight and data-efficient due to a low trainable parameter budget.
- BERT is tuned in an effective way to adapt without losing prior knowledge, while at the same time squeezing as much useful information as possible from audio-visual modalities.
- We study our model’s robustness to noise and compare its performance with a fine-tuned version and the current state-of-the-art MISA.

2. PROPOSED METHOD

Fig. 2 illustrates an overview of the system architecture. First of all, the input sequences are fed into their respective encoders to prepare for the next stage. The core component is a frozen pre-trained BERT model, which is tuned by adapter layers, without access to any other modalities. These BERT representations are combined with audio-visual information in a feedforward network (FFN) in order to perform layer-wise multimodal fusion. This process is repeated for 12 layers and the last representations are provided to a FFN to predict the sentiment score.

Frozen BERT layers: The frozen BERT model is at the core of the architecture to emphasize the importance of language.

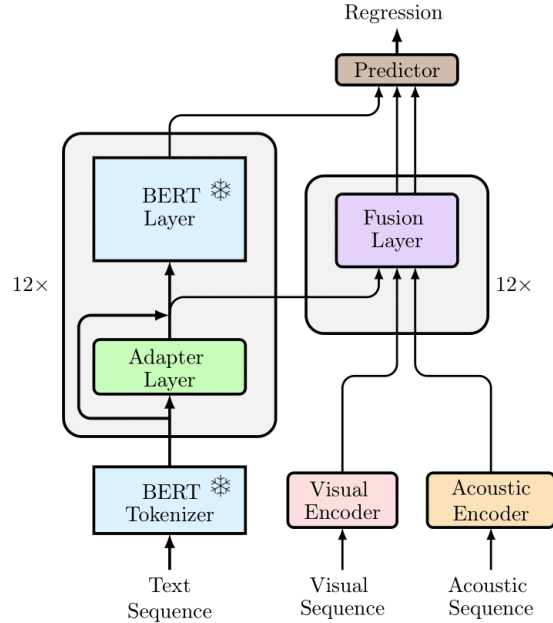


Fig. 2: Architecture of Adapted Multimodal BERT (AMB)

Both BERT tokenizer and the 12 BERT layers are kept intact during training, limiting the effects of catastrophic forgetting that can incur during fine-tuning.

Adapter layers: We use the original bottleneck adapters, introduced by Hounsby et al. [12]. Each adapter layer is composed of a linear down-projection followed by a ReLU non-linearity and then a linear up-projection to restore the original input dimensions. Residual connections are used between the input and output of each adapter layer. Instead of inserting an adapter layer both between the attention and the feedforward module, we follow [26] and only insert them after the feedforward layer. This cuts the number of additional parameters in half. Our adapter layers are only responsible for adapting to the textual inputs.

Visual and Audio Encoders: Visual and audio encoders consist of transformer encoder layers that act on each modality separately to extract information from an arbitrary sequence of features and compress it in a concatenated visual-acoustic token. This token is then prepared for the next stage of layer-wise multimodal fusion. Our encoders are closely related to the approach of [4, 5, 13], with the addition of audio.

Fusion layers: For multimodal fusion FeedForward Network Fusion (FFN-Fusion) is used in a layer-wise manner, between each BERT layer. The first BERT token (known as CLS token), which is commonly used to store a semantic summary of BERT’s hidden states [3], is projected to a lower dimension and then concatenated with the modality tokens produced by the visual and audio encoders. This tensor is then fed into FFN-Fusion to output the fused representations. Although [7] and [13] also perform layer-wise multimodal fusion, both use the result to shift BERT representations in order to generate

Models	MAE (\downarrow)	Corr (\uparrow)	Acc-7 (\uparrow)	Acc-2 (\uparrow)	F1 (\uparrow)	Trainable Parameters
MMLatch (G) [27]	0.582	0.704	52.1	82.8	82.9	2.6
MuT (G) [28]	0.580	0.703	51.8	82.5	82.3	1.8
LMF (B) [29]	0.623	0.677	50.2	82.0	82.1	1.0
TFN (B) [30]	0.593	0.700	51.8	82.5	82.3	0.6
MFM (B) [21]	0.568	0.717	51.3	84.4	84.3	1.7
ICCN (B) [23]	0.565	0.713	51.6	84.2	84.2	–
MAG-BERT* (FT) [7]	0.614	0.763	50.9	84.3	84.2	110.8
MISA (FT) [6]	0.555	0.756	52.2	85.3	85.3	47.1
AMB (Ours)	0.536	0.766	53.3	85.8	85.8	8.6

Table 1: Results on CMU-MOSEI. Models indicated with (G) use glove embeddings. Models indicated with (B) use frozen BERT embeddings, and are taken from [23]. MISA and MAG-BERT use a fine-tuned (FT) BERT for feature extraction from language. MAG-BERT* is reproduced for CMU-MOSEI by the authors of this paper. Trainable parameters are in millions.

output text. We adopt a simpler approach without shifting.

Predictor: The fused representation of the last BERT and fusion layers are concatenated and fed into a classification head, consisting of a single Feedforward layer. Minimum Absolute Error loss is used for end-to-end training of the network.

3. EXPERIMENTAL SETUP

Data: The proposed model is evaluated for sentiment analysis on CMU-MOSEI [31]. It contains 23,454 YouTube video clips of reviews on movies or other topics, where each sample is manually annotated with a sentiment score, ranging from -3 (strongly negative) to $+3$ (strongly positive). Text transcriptions are segmented into words, while visual FACET and acoustic COVAREP features are collected and aligned on these words. Standard train, development and test splits are provided. For evaluation, mean absolute error (MAE) and Pearson Correlation (Corr) between model and human predictions are used for regression, while seven-class accuracy (Acc-7), binary accuracy (Acc-2) and F1-score (F1) are used for classification.

Implementation Details: The bert-base-uncased version of BERT [2] is used for all experiments. It contains 12 transformer layers, where each token of the sentence has hidden size of 768 dimensions. The tokens are prepared for BERT with the standard tokenization procedure, while the two special tokens, [CLS] and [SEP], are added at the start and in the end of each sentence respectively. The encoders used for visual and acoustic modalities are randomly initialized transformer encoder modules with 2 layers and 1 attention head. We find that prepending a learnable [CLS] token and collecting this as a semantic summary works best. After a short hyper-parameter search in the range [128, 768] for the hidden size of the adapter layers, 384 is chosen as the optimal value. Similarly, for fusion layers 220 is chosen from [160,

Models	Corr (\uparrow)	Acc-7 (\uparrow)	Train. Params
AMB no-text	0.240	41.64	8.6
AMB text-only	0.760	52.81	8.6
MISA-Adapters	0.758	52.15	8.5
MISA	0.756	52.20	47.1
AMB-FT	0.756	51.98	47.2
AMB	0.766	53.29	8.6

Table 2: Multimodal adapters vs fine-tuning. We include experiments, where the text, or the audio-visual modalities are missing. Trainable parameters are in millions.

820] as the hidden size.

For optimization, the Adam optimizer [32] is used with learning rate $5 * 10^{-5}$. Early stopping is used with patience set to 10 epochs and dropout is set to 0.2. Training takes 20 minutes on a single GTX 1080Ti NVIDIA GPU.

4. EXPERIMENTS AND RESULTS

Comparison to state-of-the-art: The results for multimodal sentiment analysis on CMU-MOSEI are presented in Table 1. For fair comparison we only compare with methods in the literature that train in one stage, without leveraging their own, separate, pre-training stage on multimodal data. We observe that the proposed model outperforms all other methods by a significant margin across all metrics. As shown in Fig. 1, models that utilize Glove embeddings (G) [33], or frozen BERT embeddings (B), have fewer trainable parameters, sacrificing overall performance. Models that rely on fine-tuning of BERT have a significantly larger amount of trainable parameters. AMB with adapters surpasses fine-tuning based approaches on a small parameter budget.

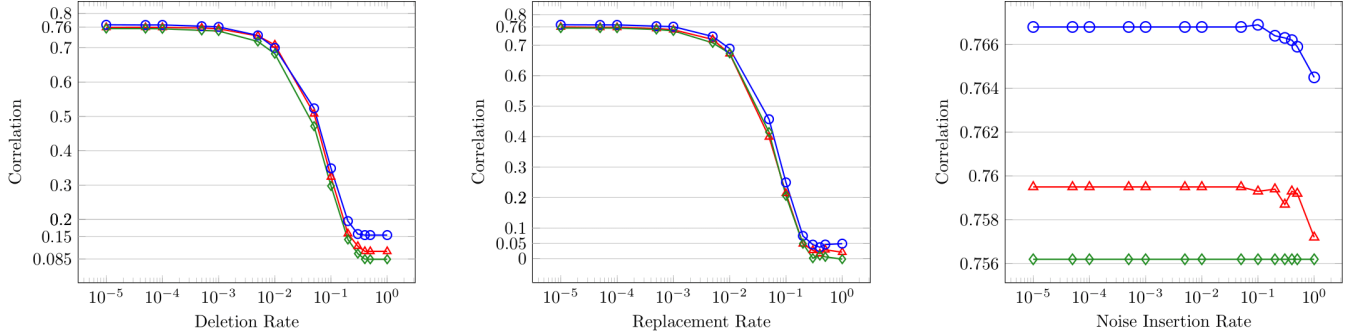


Fig. 3: Model robustness for varying levels of noise, i.e. random deletion of input tokens (left), random replacement of input tokens (middle), noise insertion to visual inputs (right). Blue \circ : AMB, Red \triangle : MISA, Green \diamond : AMB-FT

Ablation studies: Table 2 shows an ablation study on the effect of the exclusion of modalities and the effect of using adapters versus finetuning for the adaptation of the language model. Firstly, the exclusion of the textual modality significantly degrades performance for the “AMB no-text” model, which demonstrates that text is the dominant modality for this task. With the exclusion of audio-visual information in “AMB text-only” performance still declines, though to a lesser degree, indicating that the use of multimodal information is beneficial.

For the adapters versus fine-tuning experiments, an adapter based version of MISA (“MISA-Adapters”) and a fine-tuned version of AMB (“AMB-FT”) are implemented. We observe that fine-tuning is either unnecessary as in the case of MISA or even decreases model performance as in the case of AMB, revealing that some catastrophic forgetting occurs when performing fine-tuning on the text modality in this multimodal setting.

Noise Robustness: Finally, we evaluate the robustness of our model with respect to noise insertion in the visual and text modalities. When testing the robustness for the visual modality, we follow Hazarika et al. [34], who propose the insertion of multiplicative Gaussian noise to a randomly selected set of input sequence elements for a given modality. For the text modality a different approach is employed that more closely simulates real-world errors, i.e. deleting and replacing input tokens. In the token replacement experiment a percentage of input tokens is selected randomly and replaced with random tokens from the vocabulary, while for the token deletion experiment they are instead replaced with the [UNK] token. We select the best checkpoint of each model and show the average correlation over three independent runs, following [34].

Fig. 3 displays the results of the robustness tests for varying levels of input noise. The deletion, replacement and noise insertion rate refer to the probability of corrupting each element in the input sequence. When corrupting textual inputs by deleting or replacing tokens we observe that performance starts to degrade after corrupting each token with 5% probability. Steeper performance degradation occurs in the case of replacement than in the case of deletion. This sensitivity to

noise is expected, as text is the dominant modality. We observe similar robustness characteristics for AMB, MISA and AMB-FT, though adapter-based AMB appears to be somewhat more robust than its fine-tuned counterpart. In the extreme case from 50% probability and beyond AMB’s lowest point is significantly higher than the rest, verifying that it considers all modalities to make predictions. In the case of noise injection to the visual modality performance drops off for AMB and MISA at 10% noise insertion rate. We observe that noise insertion in the visual modality affects both models less than noise insertion in text. Interestingly, the AMB-FT model is not affected by visual noise, revealing that this model relies completely on text, ignoring visual cues. These results highlight that, favoring adapter-based approaches over fine-tuning when using large pre-trained language models for multimodal tasks may lead to improved model robustness and better utilization of information from less dominant modalities (that contribute less to overall performance).

5. CONCLUSIONS

In this work, AMB is proposed, a simple yet innovative model that builds on a powerful pre-trained BERT transformer encoder and avoids the pitfalls of catastrophic forgetting and modality imbalance, i.e., useful knowledge from pre-training and non-dominant modalities is leveraged effectively. Further, the use of adapters allows our model to lower the cost of trainable parameters and leads to improved robustness to various types of noise.

In the future, we plan to extend our experiments to more tasks, such as text generation from input prompts enriched with images. Moreover, exploring more sophisticated fusion methods compatible with our approach might be beneficial. The effects of shifting should also be considered. We hope that this approach will be viewed as the blueprint for designing multimodal models based on pre-trained unimodal encoders in a flexible and effective manner.

6. REFERENCES

- [1] Radford et al., “Improving language understanding by generative pre-training,” 2018.
- [2] J. Devlin, M. W. Chang, K. Lee, and K. Toutanova, “BERT: Pre-training of deep bidirectional transformers for language understanding,” in *Proc. NAACL, Vol. 1*. 2019, pp. 4171–4186, ACL.
- [3] J. Lu et al., “ViLBERT: Pretraining Task-Agnostic Visiolinguistic Representations for Vision-and-Language Tasks,” in *NeurIPS*, 2019, vol. 32.
- [4] Maria Tsimpoukelli et al., “Multimodal few-shot learning with frozen language models,” in *NeurIPS*, A. Beygelzimer et al., Eds., 2021.
- [5] C. Eichenberg et al., “MAGMA – Multimodal Augmentation of Generative Models through Adapter-based Finetuning,” *arXiv:2112.05253*, 2021.
- [6] D. Hazarika, R. Zimmermann, and S. Poria, “Misa: Modality-invariant and -specific representations for multimodal sentiment analysis,” in *Proc. 28th ACM*. 2020, MM ’20, p. 1122–1131, ACM.
- [7] W. Rahman et al., “Integrating Multimodal Information in Large Pretrained Transformers,” 2020, *arXiv:1908.05787*.
- [8] M. McCloskey and N. J. Cohen, “Catastrophic Interference in Connectionist Networks: The Sequential Learning Problem,” Academic Press, 1989, ISSN: 0079-7421.
- [9] Brown et al., “Language Models are Few-Shot Learners,” 2020, *arXiv:2005.14165 [cs]*.
- [10] B. Lester, R. Al-Rfou, and N. Constant, “The power of scale for parameter-efficient prompt tuning,” in *Proc. 2021 EMNLP*. 2021, pp. 3045–3059, ACL.
- [11] X. L. Li and P. Liang, “Prefix-tuning: Optimizing continuous prompts for generation,” in *Proc. 59th Annual Meeting of the ACL*. 2021, pp. 4582–4597, ACL.
- [12] N. Houlsby, A. Giurgiu, et al., “Parameter-efficient transfer learning for NLP,” in *Proc. 36th ICML*, 2019.
- [13] J. B. Alayrac et al., “Flamingo: a Visual Language Model for Few-Shot Learning,” Tech. Rep., 2022, *arXiv:2204.14198*.
- [14] A. Metallinou et al., “Context-sensitive learning for enhanced audiovisual emotion classification (Extended abstract),” in *2015 ACII*, 2015.
- [15] M. Wöllmer et al., “LSTM-Modeling of continuous emotions in an audiovisual affect recognition framework,” *Image and Vision Computing*, vol. 31, pp. 153–163, 2013.
- [16] A. Shenoy et al., “Multilogue-net: A context-aware RNN for multi-modal emotion detection and sentiment analysis in conversation,” in *Challenge-HML*. 2020, pp. 19–28, ACL.
- [17] S. Poria et al., “Convolutional MKL Based Multimodal Emotion Recognition and Sentiment Analysis,” in *2016 IEEE 16th ICDM*, 2016, pp. 439–448.
- [18] Y. Gu et al., “Multimodal Affective Analysis Using Hierarchical Attention Strategy with Word-Level Alignment,” in *Proc. 56th ACL*. 2018, pp. 2225–2235, ACL.
- [19] Y. Wang et al., “Words can shift: Dynamically adjusting word representations using nonverbal behaviors.,” 2019, pp. 7216–7223, AAAI.
- [20] A. Vaswani, N. Shazeer, et al., “Attention is All you Need,” in *NeurIPS*, I. Guyon et al., Eds. 2017, vol. 30, Curran Associates, Inc.
- [21] Y.-H. H. Tsai, P. Liang, et al., “Learning factorized multimodal representations,” in *ICLR*, 2019.
- [22] J. Delbrouck et al., “A transformer-based joint-encoding for emotion recognition and sentiment analysis,” in *2nd Challenge-HML*. 2020, pp. 1–7, ACL.
- [23] Z. Sun et al., “Learning relationships between text, audio, and video via deep canonical correlation for multimodal language analysis,” in *AAAI*, 2020.
- [24] W. Yu et al., “Learning modality-specific representations with self-supervised multi-task learning for multimodal sentiment analysis,” in *Proc. AAAI*. 2021, *arXiv*.
- [25] J. Kim and J. Kim, “CMSBERT-CLR: Context-driven Modality Shifting BERT with Contrastive Learning for linguistic, visual, acoustic Representations,” 2022, *arXiv:2209.07424*.
- [26] J. Pfeiffer, A. Kamath, et al., “AdapterFusion: Non-destructive task composition for transfer learning,” in *Proc. 16th ACL*. 2021, pp. 487–503, ACL.
- [27] G. Paraskevopoulos, E. Georgiou, and A. Potamianos, “Mm-latch: Bottom-up top-down fusion for multimodal sentiment analysis,” in *ICASSP IEEE*, 2022, pp. 4573–4577.
- [28] Y.H. H. Tsai, S. Bai, P. P. Liang, J. Z. Kolter, L.P. Morency, and R. Salakhutdinov, “Multimodal transformer for unaligned multimodal language sequences,” in *Proc. 57th ACL*, 2019, pp. 6558–6569.
- [29] Z. Liu et al., “Efficient low-rank multimodal fusion with modality-specific factors,” in *Proc. 56th ACL*. 2018, pp. 2247–2256, ACL.
- [30] A. Zadeh, M. Chen, et al., “Tensor fusion network for multimodal sentiment analysis,” in *EMNLP*, 2017.
- [31] A. Zadeh, P. P. Liang, et al., “Multimodal language analysis in the wild: Cmu-mosei dataset and interpretable dynamic fusion graph,” in *ACL*, 2018.
- [32] D. P. Kingma and J. Ba, “Adam: A method for stochastic optimization,” in *ICLR*, 2015.
- [33] J. Pennington, R. Socher, and C. D. Manning, “Glove: Global vectors for word representation,” in *EMNLP*, 2014, pp. 1532–1543.
- [34] D. Hazarika, Y. Li, et al., “Analyzing modality robustness in multimodal sentiment analysis,” in *NAACL*. 2022, pp. 685–696, ACL.